

1 TITLE OF THE INVENTION

2 SERVICE LEVEL AGREEMENTS AND MANAGEMENT THEREOF

3 FIELD OF THE INVENTION

4 The present invention generally relates to information technology, and more
5 particularly relates to managing one or more services for forming and complying with
6 service level agreements.

7 BACKGROUND OF THE INVENTION

8 Recent technological advances combined with market forces have resulted in the
9 creation of new services composed of other services. The term "composite service" is
10 used to describe these new services. Composite service providers pull together a plurality
11 of component services to provide a composite service. Composite services
12 conventionally span several participant organizations. Terms such as "virtual enterprise"
13 and "virtual organization" are conventionally used to describe this type of collection of
14 organizations. A participant organization may provide component services to one or more
15 virtual enterprises. Each component service provider implements a service by executing
16 a process. Implementation of a composite service requires execution of a process that
17 spans multiple organizations. The execution of such multi-organizational processes
18 conventionally involves interaction among participant organizations' computer systems.

19 By way of example and not limitation, today there are virtual enterprises reselling
20 web search services. Such virtual enterprises receive a query from a user. This query is
21 then used to query selected web search services offered by component service providers
22 of which this virtual enterprise is a client.

1 Referring to the block diagram of FIG. 1, a group of component service providers'
2 computer systems 20 comprise component service providers 20a and 20b. Component
3 service providers 20a and 20b have respective service implementations 19j and 19k.
4 Service implementations 19j and 19k may be put in communication with composite
5 service providers 10a and 10b of a group of composite service providers' computer
6 systems 10. Each composite service provider 10a and 10b may have one or more client
7 processes 13m to 13n and 13q to 13p, respectively.

8 Continuing the above-mentioned example, suppose a user of composite service
9 provider 10a places a query for a World Wide Web search. This query invokes client
10 process 13m causing a request to be sent to service implementations 19j and 19k for
11 searching the World Wide Web using respective search engines associated with these
12 services. Results from such searches may then be provided from service implementations
13 19j and 19k to client process 13m. Hence, in this example, a user executes separate
14 searches on separate search engines of separate service providers from a single query on
15 another separate service provider. In other words, a composite service provider executes
16 a business process which in turn causes component service providers to execute
17 respective business processes.

18 Accordingly, it should be understood that a component service provider may have
19 several services to offer its clients. Thus, component service providers may have a
20 platform of services available to subscribers or clients. Such services may be invoked
21 through various invocation infrastructures such as Common Object Request Broker
22 Architecture ("CORBA"), Java Remote Method Invocation ("Java RMI"), Hypertext
23 Transport Protocol ("HTTP"), among others. Moreover, this invocation may be manual;

1 for example, a phone call from a composite service provider representative to a
2 component service provider representative.

3 In the telecommunications field, Competitive Local Exchange Carriers (CLECs)
4 resell local telephone service of Incumbent Local Exchange Carriers (ILECs). Thus, a
5 CLEC may offer services of several ILECs of which it is a client and vice versa. In a
6 CLEC business model, there is interaction between ILEC and CLEC business processes.
7 By way of example and not limitation, a CLEC customer service representative may
8 interact with provisioning ILEC processes to place an order, inquire about an order, or to
9 cancel an order.

10 Accordingly, with respect to the above-mentioned Internet example and
11 telecommunications example, in order to offer their selection of services, a composite
12 service provider relies on services of its component service providers. Therefore, it is
13 incumbent upon composite service providers as clients of component service providers to
14 enter into agreements to guarantee that service needs are met. Examples of such
15 guaranteed service needs may include maximum response time and minimum throughput.
16 These agreements are referred to hereinafter as Service Level Agreements (SLAs). SLAs
17 also assist component service providers in managing their resources to meet their client's
18 needs. Without such SLAs, a component service provider may be overwhelmed by
19 requests from one client organization, which can affect service level to other clients.

20 SLAs pertain to services at an application level, as distinguished from end-to-end
21 quality of service (QoS). QoS conventionally pertains to quality parameters of a system
22 infrastructure, or more particularly network performance. A taxonomy of QoS may be

1 found in "Taxonomy of QoS Specifications," by Bikash Sabata, *et al.*, *Proceedings of*
2 *WORDS '97*, February 1997.

3 Quality objects, which are described in more detail in "Specifying and Measuring
4 Quality of Service in Distributed Object Systems," by Joseph P. Loyall, *et al.*,
5 *Proceedings of ISORC '98*, April 1998, facilitate specification monitoring of QoS
6 contracts between clients and service providers. However, this specification monitoring
7 is directed at service implementation details and not invoked functionality. Moreover, in
8 QoS contracts, a client is required to specify resource requirements. However, a client
9 may have limited knowledge of resource usage of an invoked service.

10 A QoS web server is described in "Supporting Quality of Service in HTTP
11 Servers," *Proceedings of the Seventeenth Annual SIGACT-SIGOPS Symposium on*.
12 *Principles of Distributed Computing*, June 1998. This QoS web server allows allocation
13 of server resources to specific web page requests. System capacity is represented by an
14 estimate of bytes per second served by the server. Thus, issues of guarantees to clients
15 are not addressed.

16 A product called "SilkMeter" from Segue Software, Inc. of Lexington,
17 Massachusetts, is a software system for supporting usage control in CORBA
18 environments. SilkMeter supposedly controls object usage and access based upon
19 customer-defined usage policies, and provides metering capabilities allowing software
20 owners to monitor usage activity and to bill users accordingly. However, SilkMeter does
21 not support implementation of SLAs.

22 Hewlett-Packard Company of Palo Alto, California, has announced a web QoS
23 strategy. In this announced strategy, website operators may create classes of users with

1 priorities assigned to each class, and more particularly operators may create service
2 classes and allocate capacity to each of them. However, this strategy falls short of
3 providing mechanisms that allow organizations to enter into SLAs. For example, in this
4 strategy, if two organizations are at the same priority level, then it is possible that
5 requests from only one of them will be serviced.

6 Accordingly, it would be desirable to provide specification and fulfillment thereof
7 for SLAs between organizations. Advantageously, it would be desirable for such SLA
8 specification and fulfillment to be applicable to a variety of services and implementations
9 and to facilitate deployment over existing distributed system infrastructures.

10 **SUMMARY OF THE INVENTION**

11 An aspect of the present invention is a service level agreement manager. Such a
12 service level agreement manager is disposed between one or more client process run on
13 one or more computer systems and a service implementation run on one or more other
14 computer systems. Moreover, a client process may be a service implementation. Such a
15 service level agreement manager comprises an admission controller, a performance
16 measurement module and a specification module.

17 Another aspect of the present invention is a method for service level formation.
18 More specifically, a specification module of a service level agreement manager is
19 invoked. Performance information is obtained from a performance measurement module.
20 A client provides anticipated usage information for a target service. The performance
21 information and usage information is compared to determine if a basis for forming a
22 service level agreement exists.

1 Another aspect of the present invention is a method for managing system
2 performance. More specifically, a service level agreement manager determines whether a
3 client's request is within the scope of a service level agreement. For example, it may be
4 determined whether a request is within the scope of a service level agreement in effect
5 between a requesting client and a service provider of a service implementation for which
6 this client's request is targeted. If the request is within the scope of the service level
7 agreement, the service level agreement is provided to a performance measurement
8 module and to a service organization's service implementation. Results are then obtained
9 from this service implementation in response to this request. Performance parameters
10 associated with sending a request from and receiving a response to a service level
11 agreement manager may be measured. These performance parameters may then be checked
12 against performance parameters agreed to in the service level agreement.

13 Advantageously, a service level agreement manager in accordance with the
14 present invention may be independent of service implementation with respect to
15 compatibility issues. Such a service level agreement manager need not directly monitor
16 or measure resource usage of a service provider, rather it can measure response
17 performance therefrom. Moreover, any of several well-known optimization techniques
18 can be used within such a service level agreement manager. Furthermore, such a service
19 level agreement manager may be used with any of a variety of invocation infrastructures.

20 These and other features, advantages, objects and embodiments of the present
21 invention will become more apparent from reading the following Detailed Description of
22 the Preferred Embodiments or by practicing the present invention.

1 **DESCRIPTION OF THE DRAWINGS**

2 The features of the present invention, as well as objects and advantages, will best
3 be understood from reading the appended claims, detailed description and accompanying
4 drawings where:

5 FIG. 1 is a block diagram of a group of component service providers of the prior
6 art.

7 FIGS. 2, 2A and 2B are block diagrams of exemplary embodiments of networks
8 in accordance with the present invention.

9 FIG. 3 is a flow diagram of an exemplary embodiment of SLA formation in
10 accordance with the present invention.

11 FIG. 4 is a flow diagram of an exemplary embodiment of SLA usage in
12 accordance with the present invention.

13 In the drawings, same reference numbers refer to like components throughout the
14 several figures.

15 **DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS**

16 In the following detailed description, reference is made to the accompanying
17 drawings which form a part of this detailed description, and in which, shown by way of
18 illustrative example, specific embodiments are described. These embodiments are
19 described in sufficient detail to enable those of skill in the art to practice the present
20 invention. However, it is to be understood that other embodiments of the present
21 invention not described herein in detail may be utilized. Therefore, the following
22 detailed description is not to be taken in a limiting sense.

1 Referring to FIG. 2, there is shown a block diagram of an exemplary embodiment
2 of a system 200 in accordance with the present invention. SLA manager 110 may be put
3 into or is in communication with one or more client computer systems ("clients") 100.
4 As illustratively shown, SLA manager is in communication with clients 100a and 100b.
5 Clients 100 may comprise one or more composite service providers, as described
6 elsewhere herein and may comprise one or more computer systems for running one or
7 more client processes. By communication, it is meant electrical, optical, transverse
8 electromagnetic wave, among other forms of communication.

9 SLA manager 110 may be put into or is in communication with a service provider
10 120. A service provider 120 provides service implementation 119. Service provider 120
11 is a component service provider, as described elsewhere herein, and may comprise one or
12 more computer systems for running service implementation 119.

13 Accordingly, it should be appreciated that SLA manager may contemporaneously
14 manage more than one client 100.

15 SLA manager 110 provides a front-end for service implementation 119. SLA
16 manager 110 comprises admission controller 113, performance measurement module
17 111, and specification module 112.

18 Referring to FIG. 2A, there is shown a block diagram of an exemplary
19 embodiment of a system 210 in accordance with the present invention. System 210
20 comprises clients 100a through 100d, SLA managers 110a and 110b, and service
21 providers 120a and 120b. Service providers 120a and 120b comprise respective service
22 implementations 119a and 119b. One or more invocation infrastructure 211 may be used
23 for connectivity between clients 100a through 100d and SLA manager 110a and 110b.

1 Accordingly, it should be appreciated that SLA managers 110a and 110b may be used
2 with any invocation infrastructure 211. Moreover, it should be appreciated that a client
3 100a may be able to access more than one service implementation, such as service
4 implementations 119a and 119b, by using respective SLA managers, such as SLA
5 managers 110a and 110b. Moreover, it should be appreciated that service providers 120a
6 and 120b may be a same provider.

7 Referring to FIG. 2B, there is shown a block diagram of an exemplary
8 embodiment of a system 220 in accordance with the present invention. Client 100a may
9 access one or more of service implementations 119c through 119e via respective SLA
10 managers 110c through 110e. As illustratively shown, a service implementation may be
11 coupled to a SLA manager downstream from a client and may be coupled to one or more
12 SLA managers farther downstream from the client. For example, service implementation
13 119 is couple to SLA manager 110a which is downstream from client 100a, and it is
14 coupled to SLA managers 110c and 110d which are further downstream from client 100a
15 than SLA manager 110a.

16 With continuing reference to FIG. 2, and additional reference to FIG. 3 where
17 there is shown a flow diagram of an exemplary embodiment of SLA formation in
18 accordance with the present invention, SLA formation is described.

19 At 301, a client 100a is put in communication with SLA manager 110. This
20 communication may be off-line or on-line. By off-line, it is meant a representative of a
21 client is in contact with a representative of a SLA manager, for example by calling a toll
22 free number to place an order. By on-line, it is meant that a client has contacted a SLA

1 manager using an invocation infrastructure, for example by accessing a web page for this
2 SLA manager and inputting requested information.

3 At 303 SLA specification module 112 is invoked. At 305, SLA specification
4 module 112 accesses performance information from performance measurement module
5 111. At 304, a service provider 120 presents a list of offered services or functions to a
6 client 100a, and client 100a specifies its usage parameters for each offered service it
7 selects. Examples of usage parameters include but are not limited to total number of
8 concurrent users, selected services or functions, among others. For services selected, a
9 client may specify peak invocation rate and average invocation rate. By invocation rate,
10 it is meant the number of invocations of a service per unit of time.

11 At 306, performance information obtained at 305 is compared with service(s)
12 selected and associate usage information obtained at 304 to determine if a basis for a
13 SLA exists. In this context, a basis for such a SLA is availability of resources to satisfy a
14 specified request.

15 If at 306 a basis for a SLA agreement exists, at 307 client 100a and one or more
16 service providers 120 may enter into a SLA agreement. SLA specification information
17 associated with a resulting SLA agreement is provided to admission controller 113 at
18 308.

19 If at 306 there is no basis for agreement, then a reply is sent to client 100a that
20 client provided usage parameters for identified selected services are in excess of service
21 provider's capacity.

1 With continuing reference to FIG. 2 and additional reference to FIG. 4, where
2 there is shown a flow diagram of an exemplary embodiment of SLA use in accordance
3 with the present invention, processing a service request using a SLA is described.

4 At 402, admission controller 113 determines if a request from client 100a, for
5 example, is received. If a request is received, then using an existing SLA associated with
6 this received request, admission controller 113 determines whether to accept or reject
7 such request at 403. Admission controller 113 may be configured to maximize a
8 customizable benefit function to one or more service providers 120. By way of example
9 and not limitation, this may entail allocation of resources to clients in accordance with
10 SLAs between clients and service providers. Accordingly, this decision by admission
11 controller 113 may include factors such as impact on SLAs with other clients, potential
12 benefits of servicing a request, potential penalty in rejecting a request, among others.

13 In an embodiment, a measurement and learning based implementation is used.
14 SLA manager 110 makes an initial estimate of system capacity by measuring system
15 performance under a simulated load. Thereafter, SLA manager 110, through use of
16 performance measurement module 111, continues to measure actual performance of one
17 or more service implementations to refine this initial estimate of the fraction of capacity
18 used by each function. Examples of performance measurements that may be used include
19 requests served per unit of time, bytes served per unit of time, and response time.

20 By way of example and not limitation, suppose response time is used as a
21 performance indicator. Each function f_i in the interface of a service implementation is
22 associated with a range of time. This range of time denotes minimum and maximum
23 response time for this function. An initial estimate of system capacity may be generated

1 by determining a maximum number of concurrent instances of f_i that can be executed
2 within an acceptable response time. These measurements may further be used to
3 determine the fraction of total capacity consumed by each invocation of f_i .

4 Accordingly, SLA manager 110 has opportunity to learn access patterns of its

5 clients, so an estimate, improved over that simulated by SLA manager 110, of their usage
6 variations may be expressed. SLA manager 110 can learn performance of one or more
7 service implementations under different combinations of functions invoked by clients.

8 This information may be used in combination with well-known optimization techniques

9 to improve service. Some optimization techniques that may be used are found in

10 "Reinforcement Learning for Call Admission Control and Routing in Integrated Service
11 Networks," by Peter Marbach *et al.*, in *Advances in Neural Information Processing
12 Systems*, vol. 10, the MIT Press, 1998.

13 Capacity of a service provider is denoted by a number of tokens. Each client

14 organization is assigned tokens to cover its SLA manager interaction with an associated
15 service provider. This assignment is managed within SLA manager 110, so it is
16 transparent to clients 100. A product called "Measureware" from Hewlett-Packard

17 Company of Palo Alto, California, for resource usage monitoring or a product called

18 "VAM Capacity Planner" from Zitel Corporation of Freemont, California, for capacity
19 planning, may be used to obtain an estimate for tokens needed for a request. Moreover,

20 these software tools may be used to aid in determining causes of violation of SLAs.

21 However, use of either or both of these software tools is optional.

22 At 403, admission controller 113 accepts or rejects an incoming request R_i . So

23 when a request of type R_i from client 110a is provided to an SLA manager 110,

1 admission controller 113 checks if there is a sufficient number of available tokens in
2 client 100a's account. If a sufficient number of available tokens exists in client 100a's
3 account, request R_i is accepted and the number of tokens needed for R_i is deducted from
4 client 100a's account. When request R_i is completed, this number of tokens deducted is
5 credited back to client 100a's account. However, if a sufficient number of available
6 tokens does not exist in client 100a's account at the time request R_i is received, then this
7 request is denied, and this denial is provided to client 100a at 404.

8 If request R_i is accepted at 403, then this request is provided to performance
9 measurement module 111 at 405. Performance measurement module 111 provides
10 request R_i to service implementation 119. At 406, request R_i is invoked for service
11 implementation 119. At 407, in response to this request, results are obtained from this
12 service implementation selected and provided to performance measurement module 111.
13 Performance measurement module 111 records performance measurements associated
14 with execution of this request at 408. Optionally, at 408, performance measurement
15 module 111 may further check performance measurements against SLA specification
16 requirements. At 409, results obtained in response to request R_i are provided from SLA
17 manager 110 to a client, such as client 100a, originating this request.

18 Although the present invention has been particularly shown and described with
19 respect to certain embodiments thereof, including without limitation a best mode if any, it
20 should be readily apparent to those of skill in the art that various structural, logical,
21 electrical, and other changes in form and detail may be made to these embodiments
22 without departing from the scope of the present invention as set forth in the appended

1 claims. Accordingly, the present invention is defined only by the appended claims that
2 follow this detailed description.

GOLDBECK - GOLDBECK